



Caches – basic idea

- Small, fast memory
- Stores frequently-accessed *blocks* of memory.
- When it fills up, discard some blocks and replace them with others.
- Works well if we reuse data blocks
 - Examples:
 - Incrementing a variable
 - Loops
 - Function calls



Why do caches work

- Locality principles

- Temporal locality

- Location of memory reference is likely to be the same as another recent reference.
 - Variables are reused in program
 - Loops, function calls, etc.

- Spatial locality

- Location of memory is likely to be near another recent reference
 - Matrices, arrays
 - Stack accesses



Cache performance example

- Problem (let's assume single cycle CPU)
 - 500 MHz CPU → cycle time = 2 ns
 - Instructions: arithmetic 50%, load/store 30%, branch 20%.
 - Cache: hit rate: 95%, miss penalty: 60 ns (or 30 cycles), hit time: 2 ns (or 1 cycle)
- MIPS CPI w/o cache for load/store:
 - $0.5 * 1 + 0.2 * 1 + 0.3 * 30 = 9.7$
- MIPS CPI with cache for load/store:
 - $0.5 * 1 + 0.2 * 1 + 0.3 * (.95 * 1 + 0.05 * 30) = 1.435$



Caching Vocabulary

- Miss Penalty- time to fetch a block from a lower level cache or main memory
- Block (Line) size – Amount of data in each cache address (32-256 bytes)
- Bank Size - # of sets in the cache
- Cache Size –
 - Total Data contained = (bank size) x (associativity) x (block size)
 - Usually 4-64Kb for L1, 128-512 Kb L2

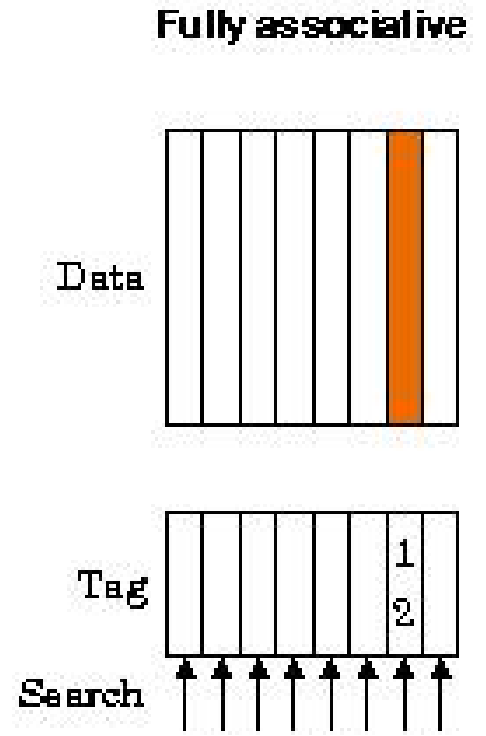
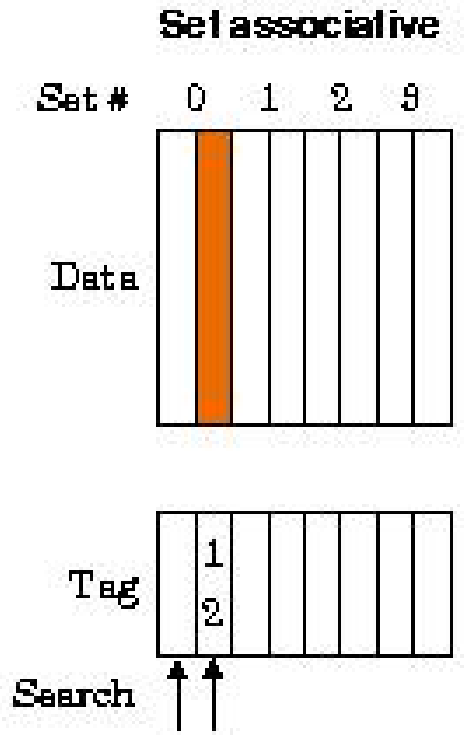
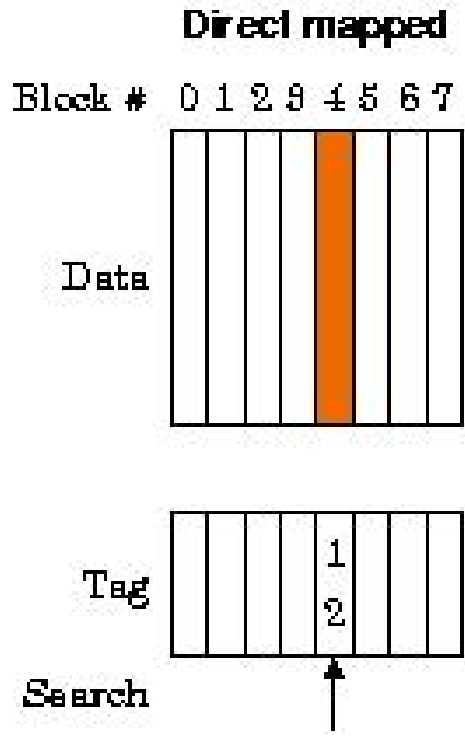


Cache types

- Direct-mapped
 - Memory location maps to single specific cache line (block)
 - What if two locations map to same line (block)?
 - Conflict, forces a miss
- Set-associative
 - Memory location maps to a *set* containing several blocks.
 - Each block still has tag and data, and sets can have 2,4,8,etc. blocks. Blocks/set = associativity
 - Why? Resolves conflicts in direct-mapped caches.
 - If two locations map to same set, one could be stored in first block of the set, and another in second block of the set.
- Fully-associative
 - Cache only has one set. All memory locations map to this set.
 - This one set has all the blocks, and a given location could be in any of these blocks
 - No conflict misses, but costly. Only used in very small caches.



More on Types





Direct-mapped cache example

- 4 KB cache, each block is 32 bytes
- How many blocks?
- How long is the index to select a block?
- How long is the offset (displacement) to select a byte in block?
- How many bits left over if we assume 32-bit address? These bits are tag bits

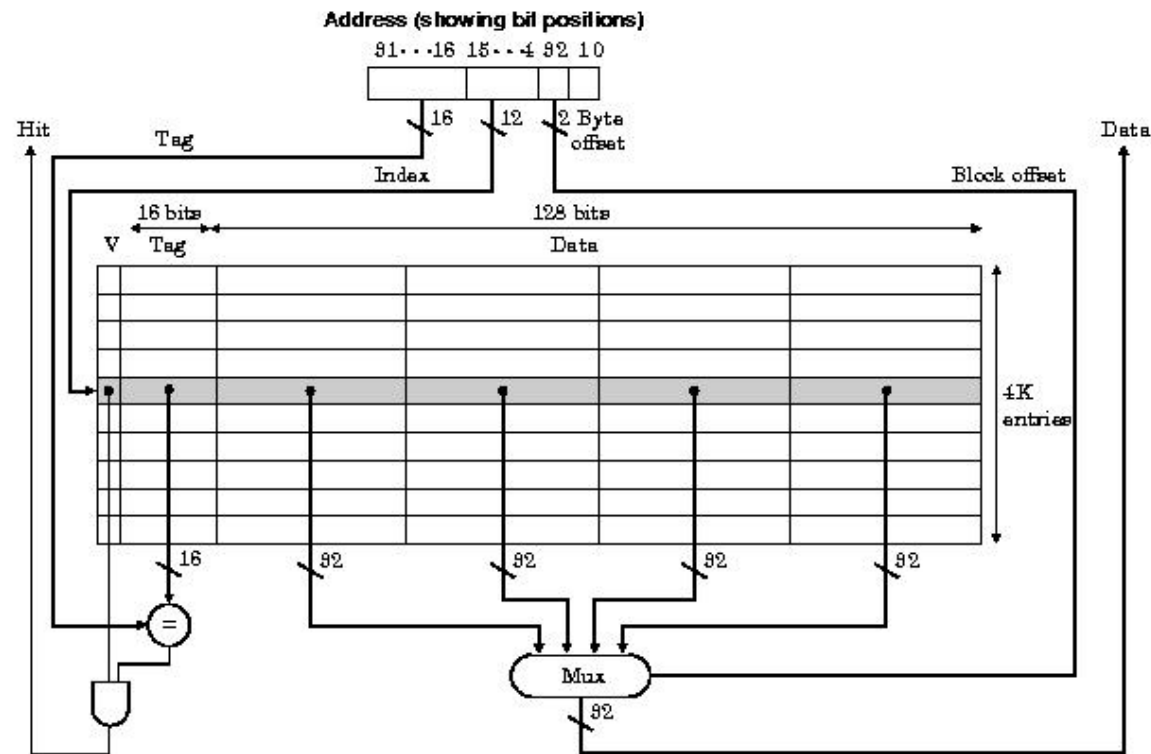


Direct-mapped cache example

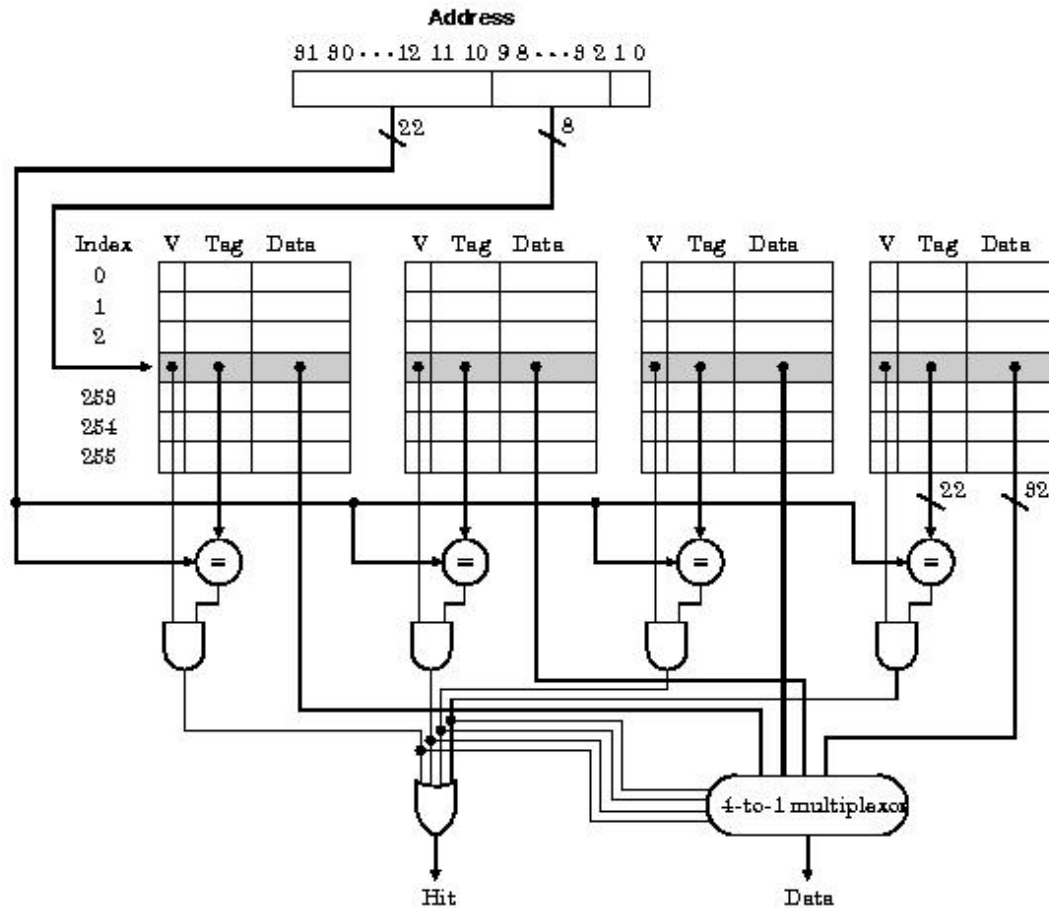
- 4 KB cache, each block is 32 bytes
 - $4 \text{ KB} = 2^{12}$, $32 = 2^5$
- How many blocks?
 - $2^{12} \text{ bytes} / 2^5 \text{ bytes in block} = 2^7 = 128 \text{ blocks}$
- How long is the index to select a block?
 - $\log_2 128 = 7 \text{ bits}$
- How long is the offset (displacement) to select a byte in block?
 - 5 bits
- How many bits left over if we assume 32-bit address? These bits are tag bits
 - $32 - 7 - 5 = 20 \text{ bits}$

Direct Mapped 4-word Block

- Address and cache:



4-way Associative 1-word block





Cache Misses: The Three Cs

- Compulsory:
 - Very first access of a block (Cold-Start Misses)
- Capacity:
 - Cache is too small to hold all blocks in the working set. Some are discarded to be retrieved later
- Conflict: (only in Direct or Set Assoc.)
 - More than n blocks map to a set in an n -way set associative cache.



Cache size

- 4 KB visible size
- Let's look at total space and overhead:
 - Each block contains:
 - 1 valid bit
 - 20-bit tag
 - 32 bytes of data = 256 bits
 - Total block (line) size: $1 + 20 + 256 = 277$ bits
 - Total cache size in hardware, including overhead storage:
 - $277 \text{ bits} * 128 \text{ blocks} = 35456 \text{ bits} = 4432 \text{ bytes} = 4.32 \text{ Kb}$
 - Overhead: 0.32 Kb (336 bytes) for valid bits and tags



Cache access examples...

- Consider a direct-mapped cache with 8 blocks and 2-byte block. Total size = $8 * 2 = 16$ bytes
- Address: 1 bit for offset/displacement, 3 bits for index, rest for tag
- Consider a stream of reads to these bytes:
 - These are byte addresses:
 - A@3, B@13, C@1, D@0, E@5, F@1, G@4, H@32, I@33, J@1
 - Corresponding block addresses $((\text{byteaddr}/2)\%8)$:
 - 1, 6, 0, 0, 2, 0, 2, 0 (16%8), 0, 0.
 - Tags: 2 for 32, 33, 0 for all others $((\text{byteaddr}/2)/8)$.
- Let's look at what this looks like. How many misses?
- What if we increase associativity to 2? Will have 4 sets, 2 blocks in each set, still 2 bytes in each block. Total size still 16 bytes. How does behavior change?...

(get notes from someone for the drawings)